## Lung Cancer Recognition Using CT Scan with CNN-VGG19 and PNN

**Mrs A S Keerthi Nayani*[1], Mrs. G.Swapnasri[2], Mr. MuddamallaNaresh[3]**
*[1]Assistant Professor,Department of ECE, Matrusri Engineering College, Hyderabad, Telangana.*
*[2]Assistant Professor,Department of ECE, Vasavi college of engineering, Hyderabad,Telangana*
*[3]Assistant Professor,Department of ECE, Matrusri Engineering College, Hyderabad, Telangana.*

**ABSTRACT**

Lung cancer is the predominant and prevalent type of cancer found in males, and it is the third most typical kind of cancer found in women. It is essential to do lung cancer screenings at primitive stage of development in order to reduce the number of deaths caused by this illness all over the world. Prior prediction of these types of diseases is very essential for better treatment analysis. The current Medical imaging techniques that are available to detect the signs of lung cancer could not recognize the disease until the severity has progressed to an advanced stage. The purpose of this research is to develop a system for the classification of lung cancer that is capable of performing diagnostic tasks in an automated fashion at the earliest stages of the disease. The examination is performed using the computed tomography (CT) imaging modalities of the lungs, and the probabilistic neural network (PNN) and VGG19 are employed for the classification task in this work. After the input lung images were preprocessed, we utilized the VGG19 classifier to identify the lung images. This classifier employs an adaptive boosting strategy that is based on the pre-trained method. The researchers found that they were able to attain accuracy rates of 96 percent and specificity rates of 97 percent in their feature selection when they used a probabilistic neural network (PNN) and VGG19-based feature selection.

**Keywords:** PNN, computed tomography, VGG19, feature selection

## I.      INTRODUCTION

Cancer is one of the major ailments where almost 75% of the population is suffering. Skin cancer [1] and Lung Cancer[2]are the mostly observed ones. Cancer of the lung is described as the uncontrolled and anomalous growing of cells that originates in one or two lungs and extends throughout the body. Cancer of the lung can begin in either of the lungs. Aberrant cells cannot proliferate in fine tissues; instead, they rapidly replicate and create tumors. Normal cells do not proliferate. In contrast to primary lung cancer, which develops in the lungs themselves, subsequent lung cancer arises somewhere else in the body and spreads from one region of the body to another until it reaches the lungs. Primary lung cancer occurs in smokers.

Men are more likely to be diagnosed with lung cancer in its primary form than women. Early symptoms of lung disease in a person's blood are typically suggestive of the beginning of the disease in the body of the patient. As the number of lung ailments in today's industrialized nations continues to rise, there has been a surge in demand for innovative techniques of accurate and early diagnosis.

According to medical professionals, smoking is the behavior that is most likely to cause lung cancer since it has an effect on the cells that make up the lungs. Smoke from cigarettes has a direct and immediate effect on the lung tissues due to the presence of harmful chemicals in it, such as carcinogens. As a consequence of this, breathing in smoke from cigarettes can have a role in the progression of lung cancer.

Primitive stages of lung cancer do typically possess no symptoms or indicators to look out for. This is because the disease has not yet fully developed. In most cases, the signs and symptoms of this condition won't start to manifest themselves until the sickness has already evolved to a more severe stage. Both the treatment and the diagnosis of lung cancer can be analyzed depending on the particular form of lung cancer, its progressed and advanced stage, and the age group of the patient. In addition, radiotherapy or chemotherapy could be utilized in the treatment of the condition. The Physical fitness & psychological strength, a number of other medical factors like development of cancer cells and willpower decide and influence patients' capacities to live longer lives. At this time, there is no specific type of diagnosis tool available in the clinical context that can be used to prevent lung cancer. However, a crucial component in the diagnosis of cancer is the earlier and more fast discovery of the disease; as a result of this, the survival rate of cancer patients will grow as a direct result of this factor.

Blood tests, radiological tests, endoscopic procedures, and biopsies are just some of the diagnostic and detection approaches that can be used for lung cancer. Biopsies can also be performed. Every kind of examination comes with its own individual set of benefits and drawbacks, in addition to some apps that are made just for it. You can get a speedy result from the test utilizing CT (Computerized Tomography), which does not cause any discomfort, and you can get information on the shape, size, and location of the tumor. An x-ray machine creates a three-dimensional image of the inside body by taking many images of the same anatomical area from a variety of angles. This image is called a computed tomography scan.

In addition to this, pathological diseases affecting the intra-thoracic region can be evaluated with the help of a CT scan. Recently

medical Experts and specialists have come up with novel technique in the process of diagnosing lung cancer. This method of diagnosis do injects a contrast medium into the circulation along with CT scans. As a direct consequence of this, the finer elements of the lung can now be seen with greater clarity. The patient's chest is scanned with a CT machine, which produces comprehensive images of the patient's chest and enables a more accurate identification of lung cancer. In this particular investigation, the CT lung Dicom (.dcm) medical input images were received from the Lung Image Database[4].The treatment of someone with lung disease is greatly aided by early detection. Because of this, testing and diagnosis are still very important to them. To begin with, radiographs of the chest (X-rays) and computed tomography (CT) scans are employed to look for potentially dangerous knobs; in either case, the possibility of polite knobs leads to erroneous conclusions. The delicate knobs are placed near to each other in the first stages. Here, a novel, profound learning-based paradigm with a variety of methodologies is required to explore the risky knobs in further detail[5].

This analysis made use of the database provided by the Consortium and (LIDC). This image dataset includes thoracic computed tomography (CT) examinations that were performed for the purposes of diagnosis and screening for cancer. Marked-up interpretations are also included (Armato et al., 2011). This database is a readily available resource that can be accessed via the internet. It serves as a global reserve for the extension, training, and approximation of computer-assisted diagnostic approaches for the disclosure and detection of cancer disease, as well as for the further therapy of patients who have cancer.
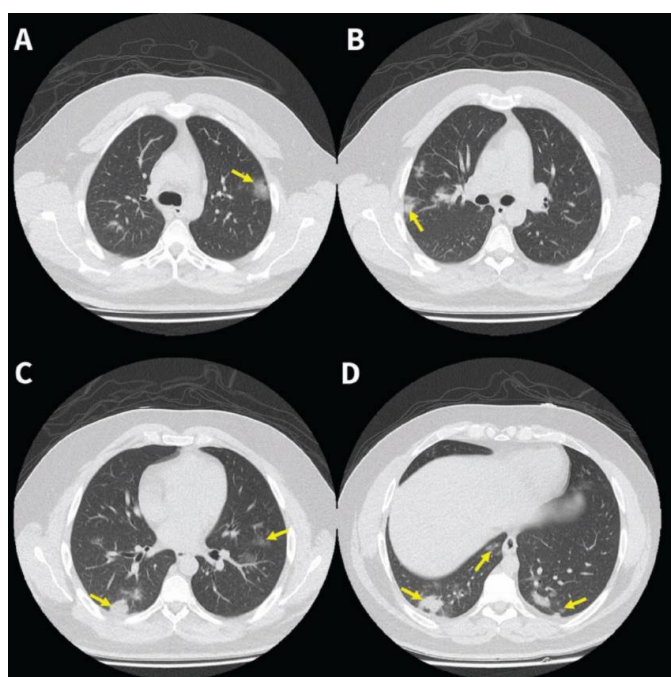


Fig.1. Lungs computerized tomography (CT) scan Images

As part of this investigation, we make an effort to employ deep learning[3] techniques to develop a method that is more efficient for the purpose of diagnosing lung cancer. The following are the most significant contributions made by the current research:

(1) Presenting an innovative approach to the treatment of lung cancer. The patient was recognized based on the images obtained from the CT scan of their lungs.

(2) Implementation of Probabilistic Neural Networks and the VGG19 model for the purpose of diagnosing lung cancer.

It is unable to generalise the model parameters across diverse lung cancer input datasets since deep learning algorithms are heuristic in nature. The input features and class labels of the various lung cancer datasets are, however, identical. The initial input vector space and the goal vector space are created by transforming the input features and the target labels. The data that passes through each layer of the deep learning[6-8] model undergoes only one simple modification[9].

Mediastinal lymph nodes are an important part of the staging of lung cancer, and the presence of metastases has a major impact on survival chances. Right now nodes are identified by a doctor but this process takes a long time and is subject to mistakes[10].Malignant nodules can be detected by recognising and measuring particular traits. The likelihood of cancer can be estimated based on the features found and their combination. Even for an experienced medical expert, this is a challenging assignment because the existence of nodules and a positive cancer diagnosis are not easily linked [11].In order to improve prediction accuracy, avoid over fitting, and reduce computational costs, dimensionality reduction may be necessary before modeling

[12]. Further, integration such systems with IoT based systems may result in better investigation even at remote locations [13].

## II.    METHODOLOGY

### 2.1 Datasets

The dataset that was used was Lung Image Database Consortium (LIDC). This type of Database Consortium archive collects CT scan images taken from actual patient cases. This Imaging Database Consortium picture array (LIDC-IDRI) consists of thoracic computed tomography (CT) scans of diagnostics as well as lung disease monitoring of labeled annotated lesions, as well as thoracic magnetic resonance imaging (MRI).CARDI is a lung cancer screening CT image database designed to aid in the development, preparation, and approximation of computer-assisted diagnostic methods for the detection and evaluation of lung cancer. CT image database consists of 1018 given data cases. Research centers &companies associated with this type of medical imaging have come up with these data cases aiding research developments. A thoracic CT clinical scan image and an associated XML file are included in each object in the collection.
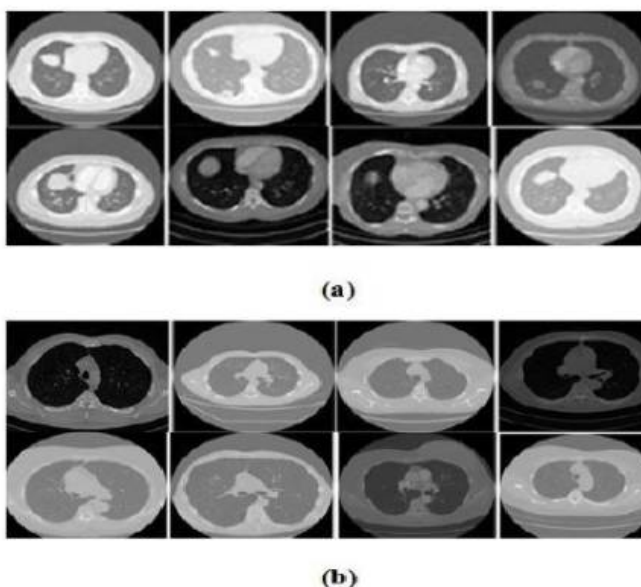


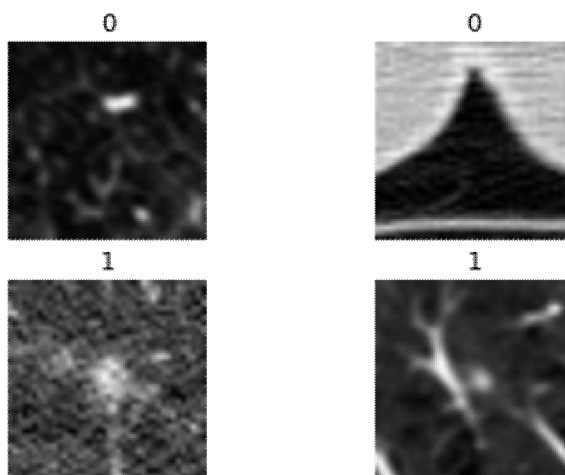Fig 2: Test images of a) Malignant b) benign from LIDC dataset



Fig 3. Dataset Sample

### 1.2. Data Preparation

Pre-processing or purification of data is a critical aspect to a Machine Learning Engineer, and the vast majority of Machine Learning Engineers put in a significant amount of effort before developing a model from the ground up. A few examples of data pre-processing techniques include outlier detection and treatment, missing value treatment, and the removal of unwanted or noisy data, to name a few. Image pre-processing, which is the same as image processing, is the term used to describe images at the lowest level of abstraction. This process does not increase the amount of image information contained in the image, but rather reduces it,

according to entropy as information metric. The goal of preprocessing is to improve the quality of the image data by suppressing unwanted distortions and enhancing some visual properties that are important for the task of subsequent processing and analysis after the image has been captured and captured. Pre-processing procedures can be divided into two categories, which are listed below. Pre-processing procedures are divided into two categories:

1.     Image Filtering and Segmentation
2.     Fourier transform and Image restoration

It is necessary to do pre-processing on input CT images in order to cut down on the amount of noise that is already there and to prepare the input pictures to be used in further processing steps such as image segmentation. As a direct consequence of this, the input images will exhibit less distortion, and the relevant characteristics of inputs will be improved. MATLAB is used to perform the preliminary processing of CT images before they are analyzed. The research takes into account cancer nodules from both the primary and secondary phases in the input database. Additionally, researchers taken into account distinct nodule types: well-circumscribed type, the juxta-pleural type, the vascularized type, and the pleural-tail type of nodules. The CT image of the patient's lung with cancer that was used as an input is shown in Figure 2.
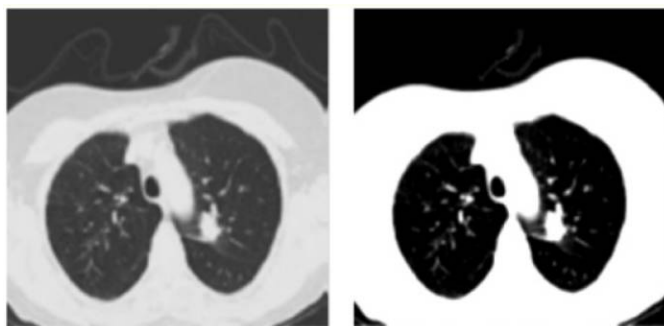


Fig 4:Binary and Cancerous Lung Image Obtained from CT

## 2.3. Modeling

**Probabilistic Neural Network**. A probabilistic neural network (PNN) is a neural network whose architecture consists of three layers referred as tri-layered, architecture being feed-forward witha one-pass training methodology. Data classification and mapping can be performed by using this type of architecture. Donald Spechthad introduced the first of this type of neural network  in the year 1990 (Specht, 1990). The fundamental idea behind probabilistic neural networks (PNN) is founded on well-known statistical principles that are obtained from Bayes' decision theory and nonparametric kernal estimators that are generated from probability density functions (pdf). While making use of Gaussian Kernel, the foremost work of the PNN lies in computing the pdf of features for each individual class given by the sample training data. In addition, following the use of the computed pdf, the Bayes decision theory is implemented in order to carry out the data classification.

It is essential to keep in mind that in PNN, weights are not "trained" in the traditional sense but rather allocated, and the weights that are now being used will not be modified in any way. During the process of training, rather than modifying the weights of preexisting vectors, new vectors are added to weight matrices. Therefore, real-time applications that are built on PNN will have a better chance of succeeding. The speed of the network is incredibly high due to the fact that matrix manipulation can be employed to carry out the running and training procedures. The PNN is able to assign the input test features to a seperate class because the class that is produced will have a greater chance of being correct than the other classes that were classed.

The most current best strategy underwent some revisions, and as a direct consequence of those revisions, a brand new idea was presented for consideration. During the pre-processing stage, the filters namely median and the Gaussian are implemented in place of the Gabor Filter.
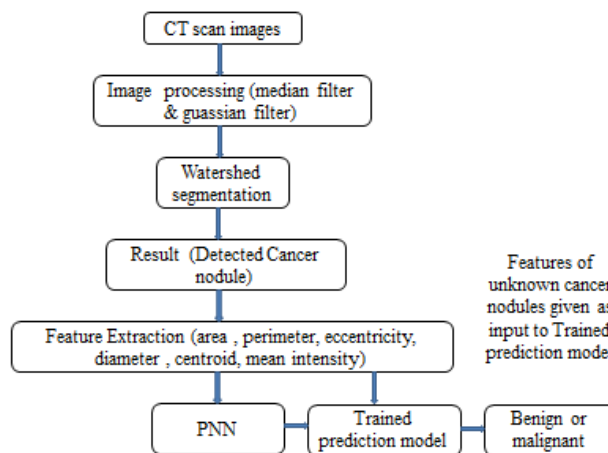
Fig 5: Probabilistic Neural Network workflow

The primary goal is to provide a visual representation of cancer nodules that have been identified. This will be accomplished in accordance with characteristics such as area, perimeter, and eccentricity. During the feature extraction process, centroid, diameter, and Pixel Mean Intensity were extracted for the confirmed cancer nodules which are used to determine their location.

Immediately following the diagnosis of the cancer nodule, the optimal model completes the process by extracting features and estimating precision. Extracted characteristics are used as training features, and as a result, a trained model is able to generate them. Then, using the learned model of prediction, the unidentified observed nodule of cancer is identified and treated accordingly. Python introduces the notions of feature recognition and feature retrieval, as well as the implementation of identification through the use of machine learning functionalities.

**VGG19**. Image recognition is performed with the help of a convolutional neural network known as VGG-19, which has 19 levels of depth. By making use of this ImageNet database it is possible to load pre-trained samples of the network. This version has trained more than one million photos in the past. The pre-trained network is able to categorize photographs into one thousand distinct object categories, some of which include animals, keyboards and mouse, pencils, and other writing implements. As a consequence of its training, the network has acquired the ability to produce rich feature representations for a wide variety of different image types. Images with a resolution of 224 pixels on each side can be uploaded to the network without issue.
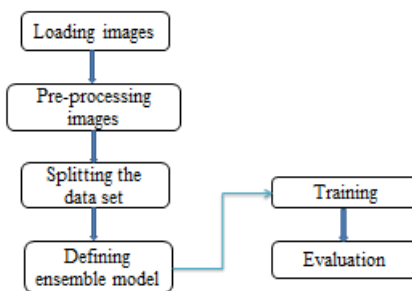


Fig 6: VGG19 workflow

The process begins with the loading of a CT scan image into CV2. Preprocessing will convert the image to grayscale, resize the resolution to 40 x 40, and convert the image into a 1D array after it has been converted. The dataset is classified into two parts: farmer one for training and later one for testing, with a 75:25 split between the two. The dataset contains a total of 2942 images, which are divided into 2206 for training and 736 for testing.

For this project, we used the VGG19 classifier, which is an adaptive boosting method based on the pretrained method. The architecture of VGG-19 is strikingly similar to that of VGG-16. The VGG-16 network has three additional convolutional layers in addition to the three already present.

After the training phase is completed, the model is evaluated against the test data in order to obtain the model's performance metrics. The model has performed admirably, with an accuracy rating of 66 percent.
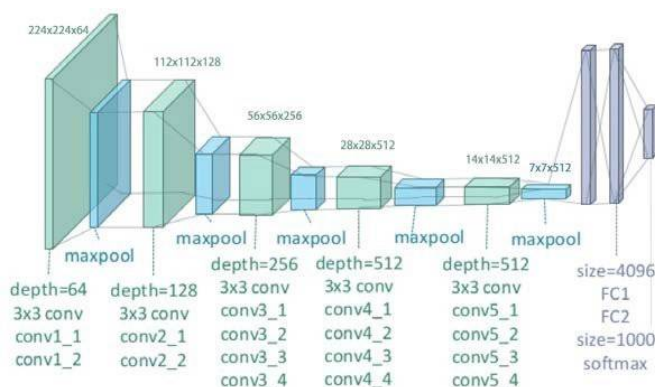
Fig 7: Network architecture of VGG-19 model

## 2.4. Validation Method

Classification report is being used to verify the model's performance, and the results are being analyzed. The precision, recall, F1, and support scores for the model are displayed in the classification report visualizer.The meaning of Precisionis the capacity of theclassifier to avoid labeling a positive instance as positive while the actual instance being negative. It is defined as the ratio of true positives to the sum of true positives and false positives for each class in the classification.Recall refers to the ability of a classifier to detect all positive instances. Defined to be same as precision for true positives but false negatives for each class in given class.The weighted harmonic mean of these two variables is what we mean when we refer to the F1 as an indicator of precision and recall. If the model's F1 score is closer to 1.0, it is anticipated that it would perform better than if it is further away from 1.0. The opposite is also believed to be true. Term "support" refers to the number of times an individual class really appears in the dataset. There were no distinguishing characteristics between the models; rather, it merely diagnoses the performance evaluation process as it takes place.

## III.    RESULTS

The performance measures are used to determine the most effective learning method (accuracy, sensitivity, and specificity). This system achieved the highest levels of performance, with average accuracy, sensitivity, and specificity rates of 96%, 97 %, and 97% respectively, in all three performance measures.

Table 1a: Probabilistic NN method Accuracy during training

| performance measures | Precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.94 | 0.99 | 0.97 | 378 |
| 1 | 0.99 | 0.93 | 0.96 | 358 |
| accuracy | | | 0.96 | 736 |
| Macro avg | 0.97 | 0.96 | 0.96 | 736 |
| Weighted avg | 0.97 | 0.96 | 0.96 | 736 |

In order to decide which learning approach is the most efficient, we employ performance measures such as accuracy, sensitivity, and specificity to evaluate each method on the best three types of image characteristics.According to the findings, VGG19 attained the best performance measures, with average rates of 66 percent accuracy, 67 percent sensitivity, and 67 percent specificity, respectively, for each of the three performance measures that were evaluated.

Table 1b: VGG-19 method Accuracy during training

| performance measures | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.74 | 0.38 | 0.50 | 378 |
| 1 | 0.57 | 0.86 | 0.68 | 358 |
| accuracy | | | 0.61 | 736 |
| Macro avg | 0.66 | 0.62 | 0.59 | 736 |
| Weighted avg | 0.66 | 0.61 | 0.59 | 736 |

## IV.     CONCLUSION AND FUTURE WORK

In the previous century, cancer rates have increased drastically due to decreased living standards, such as desk-bound lifestyles, bad dietary habits, and smoking. In wake of this, scientists and researchers have taken action to fight this terrible disease. This condition can be treated more effectively and with less risk of death if caught early enough, as per the guidelines of scientific research. An autonomous approach based on probabilistic neural networks was proposed in this paper for the best possible diagnosisof medical pictures of CT-based lungs. Because deep networks were employed to extract high-level characteristics, the suggested technique exhibited excellent classification and diagnosis accuracy.With regards to accuracy and precision, the VGG19 model is superior. Feature selection and a stacking-based strategy, which can be coupled, can be used to further increase the model's accuracy. We can also, if necessary, increase the dataset's image count to help the model perform better.

## REFERENCES

[1] R. Navid, A. Mohsen, K. Maryam et al.: Computer-aided diagnosis of skin cancer: a review, Current Medical Imaging, vol. 16, no. 7, pp. 781–793(2020).

[2] L. Hussain, W. Aziz, A. A. Alshdadi, M. S. Ahmed Nadeem, I. R. Khan, and Q.-U.-A. Chaudhry.: Analyzing the dynamics of lung cancer imaging data using refined fuzzy entropy methods by extracting different features, IEEE Access, vol. 7, pp. 64704–64721(2019).

[3] S. Lakshmanaprabu, S. N. Mohanty, K. Shankar, N. Arunkumar, and G. Ramirez.:Optimal deep learning model for classification of lung cancer on CT images,  Future Generation Computer Systems, vol. 92, pp. 374–382 (2019).

[4] Armato SG, McLennan G, Bidaut L, et al.:The lung image database consortium (LIDC) and image database resource initiative (IDRI): a completed reference database of lung nodules on CT scans. Med Phys. 2011;38:915–31.

[5] Tiwari, L., Awasthi, V., Patra, R.K., Miri, R., Raja, H., Bhaskar, N. :Lung Cancer Detection Using Deep Convolutional Neural Networks. In: Bhateja, V., Khin Wee, L., Lin, J.CW., Satapathy, S.C., Rajesh, T.M. (eds) Data Engineering and Intelligent Computing. Lecture Notes in Networks and Systems, vol 446. Springer, Singapore (2022).

[6] V. V. S. Tallapragada, D. V. Reddy, K. N. V. S. Varma and G. S. Sarma.: Improved Atrial Fibrillation Detection using CNN-LSTM ,*6th International Conference on Trends in Electronics and Informatics (ICOEI)*, 2022, pp. 1050-1055, doi: 10.1109/ICOEI53556.2022.9776825(2022).

[7] Tallapragada, V.V.S., Alivelu Manga, N., Nagabhushanam, M.V., Venkatanaresh, M..:Greek Handwritten Character Recognition Using Inception V3. In: Somani, A.K., Mundra, A., Doss, R., Bhattacharya, S. (eds) Smart Systems: Innovations in Computing. Smart Innovation, Systems and Technologies, vol 235. Springer, Singapore. https://doi.org/10.1007/978-981-16-2877-1_23(2022).

[8] Tallapragada, V. S.: Deep Learning and Its Applications in Biomedical Image Processing. In Handbook of Deep Learning in Biomedical Engineering and Health Informatics (pp. 281-302). Apple Academic Press (2021).

[9] Salama, W.M., Shokry, A. &Aly, M.H. A generalized framework for lung Cancer classification based on deep generative models. Multimed Tools Appl (2022).

[10] Gite, S., Mishra, A. &Kotecha, K. Enhanced lung image segmentation using deep learning. Neural Comput&Applic (2022).

[11] Wallis, D., Soussan, M., Lacroix, M. et al. An [18F]FDG-PET/CT deep learning method for fully automated detection of pathological mediastinal lymph nodes in lung cancer patients. Eur J Nucl Med Mol Imaging 49, 881–888 (2022).

[12] Riquelme D, Akhloufi MA. Deep Learning for Lung Cancer Nodules Detection and Classification in CT Scans. *AI*. 2020; 1(1):28-67.

[13] Tallapragada, V.V.S., Kullayamma, I., Kumar, G.V.P., Venkatanaresh, M. :Significance of Internet of Things (IoT) in Health Care with Trending Smart Application. In: Somani, A.K., Mundra, A., Doss, R., Bhattacharya, S. (eds) Smart Systems: Innovations in Computing. Smart Innovation, Systems and Technologies, vol 235. Springer, Singapore. https://doi.org/10.1007/978-981-16-2877-1_22(2022).